| | Application No. | Applicant(s) |
|---|---|---|
| ***Notice of Allowability*** | 10/020,854 | SCOTT ET AL. |
| | **Examiner** | **Art Unit** | |
| | Venkatesh Haliyur | 2619 | |

*-- The MAILING DATE of this communication appears on the cover sheet with the correspondence address--*

All claims being allowable, PROSECUTION ON THE MERITS IS (OR REMAINS) CLOSED in this application. If not included herewith (or previously mailed), a Notice of Allowance (PTOL-85) or other appropriate communication will be mailed in due course. **THIS NOTICE OF ALLOWABILITY IS NOT A GRANT OF PATENT RIGHTS.** This application is subject to withdrawal from issue at the initiative of the Office or upon petition by the applicant. See 37 CFR 1.313 and MPEP 1308.

1. ☒ This communication is responsive to *9/13/2007*.

2. ☒ The allowed claim(s) is/are *1-27*.

3. ☐ Acknowledgment is made of a claim for foreign priority under 35 U.S.C. § 119(a)-(d) or (f).

    a) ☐ All   b) ☐ Some*   c) ☐ None  of the:

        1. ☐ Certified copies of the priority documents have been received.

        2. ☐ Certified copies of the priority documents have been received in Application No. _____ .

        3. ☐ Copies of the certified copies of the priority documents have been received in this national stage application from the

           International Bureau (PCT Rule 17.2(a)).

    * Certified copies not received: _____ .

Applicant has THREE MONTHS FROM THE "MAILING DATE" of this communication to file a reply complying with the requirements noted below. Failure to timely comply will result in ABANDONMENT of this application.
**THIS THREE-MONTH PERIOD IS NOT EXTENDABLE.**

4. ☐ A SUBSTITUTE OATH OR DECLARATION must be submitted. Note the attached EXAMINER'S AMENDMENT or NOTICE OF INFORMAL PATENT APPLICATION (PTO-152) which gives reason(s) why the oath or declaration is deficient.

5. ☐ CORRECTED DRAWINGS ( as "replacement sheets") must be submitted.

    (a) ☐ including changes required by the Notice of Draftsperson's Patent Drawing Review ( PTO-948) attached

        1) ☐ hereto or 2) ☐ to Paper No./Mail Date _____ .

    (b) ☐ including changes required by the attached Examiner's Amendment / Comment or in the Office action of

        Paper No./Mail Date _____ .

    **Identifying indicia such as the application number (see 37 CFR 1.84(c)) should be written on the drawings in the front (not the back) of each sheet. Replacement sheet(s) should be labeled as such in the header according to 37 CFR 1.121(d).**

6. ☐ DEPOSIT OF and/or INFORMATION about the deposit of BIOLOGICAL MATERIAL must be submitted. Note the attached Examiner's comment regarding REQUIREMENT FOR THE DEPOSIT OF BIOLOGICAL MATERIAL.

**Attachment(s)**

1. ☐ Notice of References Cited (PTO-892)

2. ☐ Notice of Draftperson's Patent Drawing Review (PTO-948)

3. ☒ Information Disclosure Statements (PTO/SB/08), Paper No./Mail Date _____

4. ☐ Examiner's Comment Regarding Requirement for Deposit of Biological Material

5. ☐ Notice of Informal Patent Application

6. ☒ Interview Summary (PTO-413), Paper No./Mail Date *11/09/2007* .

7. ☒ Examiner's Amendment/Comment

8. ☒ Examiner's Statement of Reasons for Allowance

9. ☐ Other _____ .

## EXAMINER'S AMENDMENT

1. An examiner's amendment to the record appears below. Should the changes and/or additions be unacceptable to applicant, an amendment may be filed as provided by 37 CFR 1.312. To ensure consideration of such an amendment, it MUST be submitted no later than the payment of the issue fee.

2. Authorization for this examiner's amendment for page 1 of the specification was given in a telephone interview with the examiner by applicant's representative Thomas F. Brennan (Reg. No. 35,075) on 11/09/2007.

3. The application has been amended as follows:

Cross-References to Related Inventions

The present invention is related to the following application, which is filed on even date herewith, and which is incorporated herein by reference:

U.S. Pat. Appl. Serial No. 10/017488 (now US Pat: 6,925,547), filed December 14, 2001, entitled "REMOTE ADDRESS TRANSLATION IN A MULTIPROCESSOR SYSTEM" (Attorney Docket No. 499.709US 1).

**Allowable Subject Matter**

4.    The following is an examiner's statement of reasons for allowance:

Claims 1-27 are allowed over prior art.

The prior art of record fails to teach and render obvious the limitations as

in the independent claims 1,12,22,23 and dependent claims for a system and

method related to Node Translation and Protection in a Clustered Multiprocessor

System.

Some multiprocessor systems employ block transfer engines to transfer

blocks of data from one area of memory to another area of memory. Block

transfer engines provide several advantages, such as asynchronous operation

and faster transfer performance than could be achieved by the processor.

Unfortunately, existing block transfer engines suffer from problems that limit their

utility. For example, since address translations are performed in on-chip TLBs at

the requesting processors, external block transfer engines are prevented from

being programmed using virtual addresses. Instead, with existing block transfer

engines, user software makes an operating system (OS) call to inform the OS

that it wants to transfer a particular length of data from a particular source

(specified by its virtual address) to a particular destination (also specified by its

virtual address). In response, the OS first checks whether it has address

translations for all of the virtual addresses, and then generates separate block-transfer requests for each physical page. For example, if the virtual address range spans 15 physical pages, an OS may have to generate separate queued block-transfer requests to cause 15 separate physical transfers. The large amount of overhead associated with such OS intervention means that much of the advantage that is associated with performing the block transfer in the first place is lost.

Clustered multiprocessor systems include collections of processing machines, with each processing machine including a single processor system or distributed shared memory multiprocessor system. Clustering advantageously limits the scaling required of a single OS, and provides fault containment if one of the machines should suffer a hardware or OS error. In a clustered system, however, memory accesses to remote machines are typically performed via a network interface I/O device that requires OS intervention to send messages, and can target only specific memory buffers that were reserved for this communication at the remote machine. Thus, memory must be specifically "registered" by a user process on the remote machine, which prevents the memory on the remote machine from being accessed arbitrarily. Also, state must be set up on the remote machine to direct the incoming data, or the OS on the remote machine must intervene to handle the data, copying the data at least once. More recently, some network interface cards have been designed to support user-level communication using the VIA, ST or similar "OS bypass"

interface. Such approaches, while successful in avoiding OS intervention on
communication events, do not unify local and remote memory accesses. Thus,
programs must use different access mechanisms for intra-machine and inter-
machine communication.

Thus, there is a need for a node translation mechanism for communicating
over virtual channels in a clustered system that supports user-level
communications without the need for OS intervention on communication events.
There is also a need for a node translation mechanism that unifies local and
remote memory accesses, thus allowing user programs to use the same access
mechanisms for both intra-machine and inter-machine communications. Such a
mechanism would allow communication with other nodes in a local machine to be
handled in the same way as communications with nodes in remote machines.
There is also a need for a node translation mechanism which supports low
overhead communications in scalable, distributed memory applications that
seamlessly span machine boundaries, provides protection, and supports remote
address translation.

The invention in the instant application overcomes the above-identified
problems as well as other shortcomings and deficiencies of existing technologies
by providing a method of node translation for communicating over virtual
channels in a clustered multiprocessor system using connection descriptors
(CDs). The system includes local and remote processing element nodes and a
network interconnect there between for sending communications. The method

includes assigning a CD to a virtual connection (the CD is a handle specifying an

endpoint node for the virtual connection), defining a local connection table (LCT)

to be accessed using the CD to produce a system node identifier (SNID) for the

endpoint node, generating a communication request including the CD, accessing

the LCT using the CD of that communication request to produce the SNID for the

endpoint node for the connection in response to the communication request, and

sending a memory request to the endpoint node. The memory request is sent to

the local processing element node if the endpoint node is the local processing

element node, and is sent over the network interconnect to the remote

processing element node if the endpoint node is the remote processing element

node.

Another aspect of the invention relates to a node translation apparatus for

a clustered multiprocessor system, including a memory and communication

engine (CE). The memory stores a local connection table (LCT) having a plurality

of entries indexed by a connection descriptor (CD), each entry of the LCT storing

a system node identifier (SNID) for the endpoint of a virtual connection. The CE

receives a communication request including a CD from a user process, accesses

the LCT using the CD of the communication request to produce the SNID for the

endpoint node for the virtual connection, and sends a memory request to the

endpoint node identified using the LCT. The memory request is sent internally to

the endpoint node if the endpoint node is located within the local processing

element node, and is sent over a network interconnect to the endpoint node if the endpoint node is located within the remote processing element node.

5.     Any inquiry concerning this communication or earlier communications from the examiner should be directed to Venkatesh Haliyur whose telephone number is 571-272-8616. The examiner can normally be reached on Monday thru Friday 8:30AM to 4:30PM. If attempts to reach the examiner by telephone are unsuccessful, the examiner's supervisor, Edan Orgad can be reached on 571-272-7884. The fax phone number for the organization where this application or proceeding is assigned is 571-273-8300.

6.     Information regarding the status of an application may be obtained from the Patent Application Information Retrieval (PAIR) system. Status information for published applications may be obtained from either Private PAIR or Public PAIR. Status information for unpublished applications is available through Private PAIR only. For more information about the PAIR system, see http://pair-direct.uspto.gov. Should you have questions on access to the Private PAIR system, contact the Electronic Business Center (EBC) at 866-217-9197 (toll-free).

Venkatesh Haliyur

Patent Examiner

EDAN . ORGAD
SUPERVISORY PATENT EXAMINER